

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)**ScienceDirect**

Procedia Computer Science 89 (2016) 209 – 212

**Procedia**  
Computer Science

Twelfth International Multi-Conference on Information Processing-2016 (IMCIP-2016)

## Point Cloud Mapping Measurements using Kinect RGB-D Sensor and Kinect Fusion for Visual Odometry

N. Namitha<sup>a</sup>, S. M. Vaitheeswaran<sup>b,\*</sup>, V. K. Jayasree<sup>a</sup> and M. K. Bharat<sup>b</sup><sup>a</sup> College of Engineering, Alappuzha Dt., Kerala State 688 541, India<sup>b</sup> National Aerospace Laboratories, Bengaluru, India

### Abstract

RGB-D camera like Kinect make available RGB Images along with per-pixel depth information in real time. This paper uses the Kinect Fusion developed by Microsoft Research for the 3D reconstruction of the scene in real time using the MicroKinect Camera and applies it as an aid for Visual Odometry of a Robotic Vehicle where no external reference like GPS is available.

© 2016 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license

(<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Peer-review under responsibility of organizing committee of the Organizing Committee of IMCIP-2016

**Keywords:** Kinect; Kinect Fusion; Odometry; Robotics; Vision.

### 1. Introduction

Visual Motion Research is at the heart of Computer Vision techniques for advancing a variety of critical technologies. These include three dimensional reconstruction, tracking, surveillance, recognition, navigation and control to name a few. In the domain of reconstruction methods, using structure from motion techniques, optical flow<sup>1</sup> provides information about the motion of a three dimensional scene using two dimensional projections. For a single view camera using optical flow limitation however exists, preventing the reliable estimation of motion of a three dimensional scene. The solution to this is provided by using stereo or multiple cameras which concurrently estimate both structure and motion increasing the robustness of these techniques. The downside of the methods is the complexity, requiring mapping over a number of frames, iterative refinement steps such as bundle adjustment, reduction of uncertainty, lack of smoothness due to 2D parameterization from surfaces and so on.

With the availability of time of flight RGB-D (colour-depth) cameras using structured light sensing, in recent years, access to three dimensional information in real time and frame rates has become possible. This warrants a relook in the reconstruction methodologies and formulations. Since a depth camera is available, the sensor provides structure information and surface estimation are not needed. The earlier methods were based on the displacements of the colour pattern. There are several other advantages, one of these is the introduction of the Kinetic Fusion Algorithm<sup>2</sup> and its extensions that enable 3D dense scanning using a moving volume and the representation to an octree which does not use the standard Iterative Closest Point (ICP) making it possible to work in real time.

\*Corresponding author. Tel.: 918025086704; Fax: 918025268546.

E-mail address: [smvairu@nal.res.in](mailto:smvairu@nal.res.in)

In this paper, we propose to use the Kinect Fusion developed by Microsoft Research for the 3D reconstruction of the scene in real time using the Micro Kinect Camera and apply it as an aid for Visual Navigation of a Robotic Vehicle where no external reference like GPS is available.

## 2. Related Work

Vision sensors have shown potential to provide pose information (through dead reckoning) in structured areas and in cluttered environments where the traditional GPS can get degraded and/or get denied. Among the different approaches three common methods have received increased attention: Vision Odometry<sup>3</sup>, Vision Based Simultaneous Location and Mapping<sup>4</sup> and Structure from Motion<sup>6</sup> applications. Among these techniques used, Vision Odometry is a low latency and low cost approach and outperforms the other two approaches in terms of computational complexity and hardware. In contrast, the latter two methods are computationally intensive and require mapping over a number of frames, iterative refinement steps such as bundle adjustment, reduction of uncertainty and so on.

Visual Odometry can be split into two categories: a) Sparse Odometry and b) Dense Odometry. The sparse methods extract a set of sparse points using feature detectors like Harris, FAST or feature descriptors such as Speeded up Robust features (SURF), Scale Invariant Feature Transform (SIFT). Correspondences between the features are established between successive frames in time. A match is assumed if the error between the patch of points is minimal. From a set of good correspondences the Transformation matrix describing the set of translations and rotations are computed, frame to frame.

In contrast, the Dense Odometry approaches use a dense set of data or the whole image data. The Transformation data is obtained from the photometrical error between the frames. Alternately, in a different approach, the geometrical error between the surfaces is obtained to describe the rigid body transformations. A drawback of this approach is that they need structured surfaces and a further problem is that they require a computationally expensive nearest neighbour search to create point correspondences. To overcome this issue, the 3D surfaces are represented as 2D depth maps. The correspondence for one point in a second depth map is found by applying the rigid body motion and projecting it to 2D coordinates. These methods suffer from long term drift and consequently accumulate errors in the motion estimation process.

Fortunately the availability of Microsoft Kinect and other RGB-D sensors using structured light sensing has resulted access to three dimensional information leading to more reliable information about the structure reconstruction process and at in real time. This has opened up new possibilities and opportunities in the field of tracking and navigation. One of the significant advances is the development of the Kinetic Fusion Algorithm<sup>5</sup> which has demonstrated its potential for real time dense scanning of indoor static scenes. A further gain has been the representation of the algorithm in terms of hierarchical octrees (pyramidal structure) to achieve faster computations and parallelization. Taking advantage of this, the present paper proposes to use a RGB-D camera, readily available for visual odometry application. Since the structure information is readily available, we compute the camera motion from two consecutive frames.

## 3. Methodology

The present work is based on the Kinect Fusion to generate a 3D robust reconstruction of the environment of a robot in real-time. This is obtained by moving the Microsoft Kinect sensor around the real scene. The input to the Kinect Fusion algorithm is a temporal sequence depth maps from the Kinect sensor. This algorithm uses only the depth maps and no color information, and do not interfere with lighting conditions, allowing Kinect Fusion to function even in complete darkness. The following Fig. 1 shows the obtained Kinect Fusion result.

The algorithm runs at 30 fps in real time and provides a surface representation for each current depth frame refining a global model

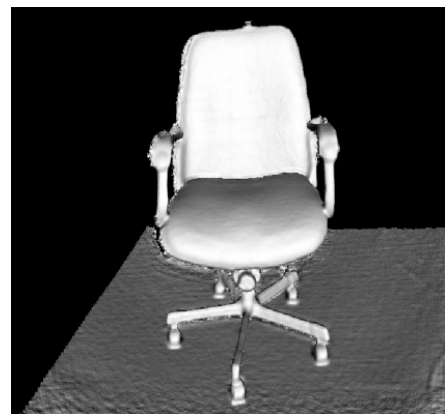


Fig. 1. A chair reconstructed with Kinect Fusion.

by merging the new surface at each time step using the ICP algorithm with it. A bilateral filter is used to smooth the image and remove noise from the depth map by replacing the intensity value at each pixel in the image with a weighted average of intensity values from nearby pixels. The result is a smoother depth map that preserves sharp edges. The depth maps are converted to a 3D point cloud with vertex and normal information using a hierarchical octree pyramid algorithm so that the result contains different levels of details.

For matching the consecutive point clouds into a single 3D model, the 6DOF camera pose estimates are obtained at each time  $t$  during the mapping process. For aligning the consecutive point clouds (we can call it as source and target point clouds) we have to choose the first point cloud as the reference model. Then the new arriving point clouds are aligned with this model using an iterative method called ICP. The Iterative Closest Point (ICP)<sup>7</sup> algorithm is used to minimize the Euclidian squared errors between the points of the source and target acquired point clouds. This allows the algorithm to run faster when it is being parallelized on GPU.

After this first match, the modified ICP algorithm computes the error between the two point clouds using a point-to-plane metric method instead of the point-to-point standard one. In this ICP variant, the sum of the squared distances between each point of the source cloud and the tangent plane at the corresponding destination point on the target cloud is minimized over the course of several iterations, until a close enough match has been found.

Consecutive point clouds are back projected onto the camera image frame of the model. Points falling on the same pixel are considered matched. The Kinect Fusion uses the a point-to-plane metric method<sup>8</sup> instead of the point-to-point standard<sup>9</sup> wherein the sum of the squared distances between each point of the source cloud and the tangent plane at the corresponding destination point on the target cloud is minimized over the course of several iterations, until a close enough match is found.

#### 4. Odometry Registration

In this case, we synchronize the acquisition of the point clouds data obtained at each time step ' $t$ ' with the odometry transformation matrix ' $T$ ' comprising a set of rotations and translations. In small angle approximate form the matrix is written as

$$T = \begin{pmatrix} 1 & \gamma & \beta & t_x \\ \gamma & 1 & -\alpha & t_y \\ -\beta & \alpha & 1 & t_z \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

that moves and rotates the points  $s_i$  such that the distance between the moved points and the tangential planes defined by the surface normals  $n_i$  is minimal:

$$\operatorname{argmin}_T \sum_i^N (Ts_i - d)_i \cdot n_i.$$

Here, is the scalar product, and  $n_i$  is the normal vector on  $d_i$ . Surface normals are readily available from Kinect Fusion depth images.

#### 5. Experiments and Results

To decide about the quality of every generated point cloud, experiments is carried out first on synthetic data sets and an open source software "Cloudcompare"<sup>10</sup> is used to compute a score that reflects how much the generated 3D point cloud corresponds to the ground truth 3D model. Figure 1 shows a typical data cloud obtained from Kinect Fusion. The mean distance in meters between the present approach and the reference model is 0.038 meters.

Figure 2 shows the RGB images and the depth images obtained from Kinect which is used to reconstruct the 3D point cloud map in real time. Cloud cap mean distances for each of the frame grab clouds are 0.0824 m, 0.087 m and 0.073 m respectively. Figure 3 shows the 6DOF Trajectory obtained by inversion of the point cloud Transformation matrix of the real time 3D point clouds. The result obtained had good agreement with the truth model.



Fig. 2. (a) RGB Image; (b) Depth Map; (c) Point Cloud from Kinect.

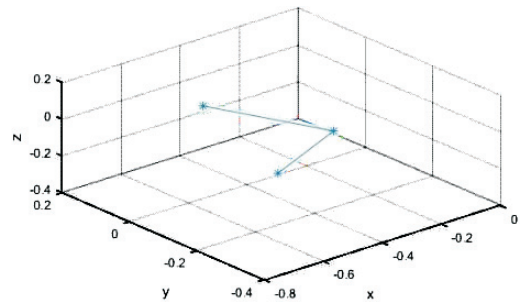


Fig. 3. 6DOF Trajectory for real time experiments in Fig. 2.

## Notes

1. See <http://www.xbox.com/en-US/kinect> and <http://www.primesense.com/>
2. See <http://www.canesta.com/> and <http://www.mesa-imaging.ch/>
3. See <http://www.primesense.com/>.
4. See <http://www.xbox.com/en-US/kinect>

## References

- [1] J. Shi and C. Tomasi, Good Features to Track, *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 593–600, (1994).
- [2] Richard A. Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J. Davison, Pushmeet Kohli, Jamie Shotton, Steve Hodges and Andrew Fitzgibbon, Kinect Fusion: Real-Time Dense Surface Mapping and Tracking, *IEEE ISMAR*, October (2011).
- [3] D. Nistér, O. Naroditsky, J. Bergen, Visual Odometry, Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004, CVPR 2004, vol. 1, pp. 1–652, 27 June (2004).
- [4] J. Artieda, J. M. Sebastian, P. Campoy, J. F. Correa, I. F. Mondragón, C. Martínez and M. Olivares, Visual 3-D SLAM from UAVs, *Journal of Intelligent and Robotic Systems*, vol. 55(4–5), pp. 299–321, 1 August (2009).
- [5] Z. Zhang, Microsoft Kinect Sensor and its Effect, *MultiMedia, IEEE*, vol. 19(2), pp. 4–10, February (2012).
- [6] J. Oliensis, A critique of Structure-from-Motion Algorithms, *Computer Vision and Image Understanding*, vol. 80(2), pp. 172–214, 30 November 2000.
- [7] E. Ezra, M. Sharir and A. Efrat, On the Performance of the ICP Algorithm, *Computational Geometry*, vol. 41(1), pp. 77–93, 31 October (2008).
- [8] S. Y. Park and M. Subbarao, An Accurate and Fast Point-to-Plane Registration Technique, *Pattern Recognition Letters*, vol. 24(16), pp. 2967–76, 31 December (2003).
- [9] F. Bosché, Automated Recognition of 3D CAD Model Objects in Laser Scans and Calculation of As-Built Dimensions for Dimensional Compliance Control in Construction, *Advanced Engineering Informatics*, vol. 24(1), 107–18, 31 January (2010).
- [10] D. Girardeau-Montaut, Cloud Compare-Open Source Project, OpenSource Project (2011).